

A priori analysis: an application to the estimate of the uncertainty in course grades

G.L. Lippi^{1,2}

¹ *Institut Non Linéaire de Nice,
Université de Nice Sophia Antipolis*

² *CNRS, UMR 7335
1361 Route des Lucioles, F-06560 Valbonne, France
Gian-Luca.Lippi@inln.cnrs.fr
(Dated: March 4, 2014)*

The *a priori analysis* (*APA*) is discussed as a tool to assess the reliability of grades in standard curricular courses. This unusual, but striking application is presented when teaching the section on data treatment of a Laboratory Course to illustrate the characteristics of the *APA* and its potential for widespread use, beyond the traditional Physics Curriculum. The conditions necessary for this kind of analysis are discussed, the general framework is set out and a specific example is given to illustrate its various aspects. Students are often struck by this unusual application and are more apt to remember the *APA*. Instructors may also benefit from some of the gathered information, as discussed in the paper.

I. INTRODUCTION

Teaching statistical data treatment techniques, generally done within the framework of a laboratory course in the Physics Curriculum, is a task that is known to many of us as an unrewarding one. On the one hand, the topics to be treated need attention to detail and to the basic assumptions which ensure the validity of the whole analysis. On the other hand, students often find the topic boring and unappealing. Nonetheless, it is a job that needs to be accomplished, in the same way one has to learn many other basic techniques indispensable for the conduct of one's work. Thus, it is always helpful to try and find ways of rendering the material more appealing to students, for example by using some unusual and unexpected applications of data treatment.

One such example, which infallibly attracts students' attention, is the use of the *A Priori Analysis* (*APA*) to estimate the uncertainty on the grade which they receive in the same laboratory course where I present this material. Besides attracting the students' attention, this example never fails to arouse some curiosity – often disguised – since the concept of a grade not being given with perfect certainty appears to surprise a number of students (and perhaps unsettle some of them a little bit).

Thus, in addition to giving students an illustration of the *APA* which brings the message very close to home, this application carries a triple weight: 1. giving a practical implementation of a concept which may otherwise be relegated to the category of techniques to be set aside (and forgot)[3]; 2. introducing the pedagogically important concept that grades, like everything else in real life, are affected by an intrinsic uncertainty; 3. showing that the tools that are learnt within the Physics Curriculum have an application to real life.

Instructors may also find that this application of the *APA* offers some useful information, as illustrated in the course of the paper and discussed in more detail in the conclusions.

II. A PRIORI ANALYSIS

Although not always included among the experimental analysis techniques taught in the basic curriculum, the *APA* is a powerful method which makes it possible to identify the different sources of uncertainty and quantify their influence on the outcome of an experiment. Two applications of the *APA* are obvious: 1. estimating the size of errors which can be expected **before** performing an experiment – particularly interesting for long and complex or expensive experiments – and 2. evaluating the influence of the individual error sources, thus permitting the identification and/or removal of the strongest uncertainty contributions. In this sense in Physics one could say that the *APA* is most useful in experiment design and/or performance evaluation. [4] However, in addition to these two more immediate applications, the *APA* becomes a valuable tool of error assessment whenever the *a posteriori* analysis – through repetition of the experiment – is not possible[5].

Since the concept of repetition may be arguable in the example that I will discuss (tests can be repeated, and multiple tests are regularly administered in a course), it is useful to recall the constraints which need to be fulfilled for the *a posteriori* analysis to hold [1]. Indeed, since mean and standard deviations are estimated through the repetition of the measurement, *each outcome must be statistically independent from the previous ones* (i.e., the fluctuations from one measurement to the other be random). If this condition is not fulfilled, the estimated mean and standard deviation do not hold and the results lose significance.

The statistical independence hypothesis certainly does not hold for at least two reasons:

1. it is virtually impossible to devise tests which are perfectly equivalent, thus the *experimental conditions are not constant*;
2. one cannot repeat the *experiment* by having students take repeated (equivalent) tests, since the

(desirable, and normally observed) progressive improvement accompanying successive tests would skew the results (particularly the standard deviation). Hence, the condition of statistical independence is clearly violated.

While one could think of finding ways of compensating for the difficulty presented in (1.) by taking large ensemble averages (i.e., repeated tests for each student) which could smooth out the differences if the tests are sufficiently well constructed, the obstacle presented by point (2.) is unsurmountable and even contradicts the possible solution just envisaged for point (1.). Indeed, in addition to obtaining meaningless results with test repetition, one cannot even think of replacing multiple tests with ensemble averages – on the class size – taken on a single test, since their indicators (average and standard deviation) cannot give any information on the grade of each individual student! The variability in level for an entire class is typically **much** wider than the accuracy with which we can estimate the grade for an individual, since the former represents the excursion in achievement due to different levels of individual ability, assiduity, performance, engagement, etc. Indeed, if this were not the case individual grades would be meaningless.

The *a priori* estimate of the uncertainty is therefore an interesting indicator which, far from providing the correct uncertainty, gives for it at least a reasonable estimate. As always true for *a priori* uncertainties, the quality of the final outcome – at the end of the analysis – is strictly related to the reliability of the guesses for the uncertainty of each test component. The appropriation of this concept is pedagogically very important, as it teaches the student to critically analyze the problem and shows that within the framework of an *APA* a critical eye and repeated tests (with different estimates for the different error components) play a major role in the process. Testing various estimated initial uncertainties is all the more useful the more numerous the tests which compose the final grade (although the case I discuss in detail turns out to be very simple). As such, lab courses may be the most interesting examples, but the technique is applicable to any kind of course.

III. MATHEMATICAL FORMULATION

One can generally formulate the problem as follows. An ensemble of N tests of different nature – lab report, written test, final exam, etc. –, with individual grades G_j , combine with weight coefficients w_j to give the global grade G :

$$G = \sum_{j=1}^N w_j G_j. \quad (1)$$

Each category of test may itself be subdivided into different subtests and result therefore from an average

over homogeneous grades:

$$G_j = \sum_{k=1}^M \frac{G_{j,k}}{M}, \quad (2)$$

where M is the number of homogeneous tests in each category. A concrete example[6] can illustrate the structure of the grades more easily. Suppose that the grade be composed of:

1. Work performance during the lab session;
2. Quality of the reporting in the labbook;
3. Evaluation of the lab report;
4. Final exam;

thus of $N = 4$ different kinds of tests. Each category of test may contain repetitions of individual tests. For instance, a student will do several, M , experiments and therefore will have at least M notes in category 1.

We first define the uncertainty in the homogeneous test category as

$$\sigma_{G_j} = \frac{1}{M} \sqrt{\sum_{k=1}^M \sigma_{j,k}^2}, \quad (3)$$

where $\sigma_{j,k}$ represents the uncertainty estimated (in the *APA*) for each individual test within a category. Thus, equation 3 provides the general expression for the *a priori* estimate of the uncertainty in each grade category[7].

In most cases, however, the uncertainty can be estimated to be the same for all k tests of a certain category j (e.g., homework, lab report, etc.). In such a case the uncertainty, equation 3, simplifies to become [1]

$$\sigma_{G_j} = \frac{1}{M} \sqrt{\sum_{k=1}^M \sigma_j^2}, \quad (4)$$

$$= \frac{\sigma_j}{\sqrt{M}}. \quad (5)$$

In order to obtain the uncertainty on the final grade, it suffices to propagate the individual uncertainties σ_{G_j} through the general definition, equation 1, to obtain [1]

$$\sigma_G = \sqrt{\sum_{j=1}^N \left(\frac{\partial G}{\partial G_j} \right)^2 \sigma_{G_j}^2}, \quad (6)$$

$$= \sqrt{\sum_{j=1}^N \left(w_j^2 \frac{\sigma_j^2}{M} \right)}, \quad (7)$$

where the former expression is general and the latter applies to the case of equal estimated uncertainties within a test category (cf. equation 4).

TABLE I: The *Repetitions* column corresponds to the number of tests in the corresponding category. Different lab sessions (4 in this example) lead to one report, thus the number of lab reports is four times smaller than the number of grades in the participation or labbook sections.

Kind of test	Kind of evaluation	label (j)	w_j	Repetitions
Participation	individual	p	0.1	12
Labbook	collective	l	0.25	12
Report	collective	r	0.25	3
Oral	individual	o	0.4	1

IV. ESTIMATING THE A PRIORI UNCERTAINTIES

The most interesting, and challenging, part of the work comes when one has to determine reasonable estimates for the uncertainties to be attributed to each individual type of test[8]. For clarity, I will proceed with a concrete example: the laboratory course in which this material is presented. The structure of the course is such that students are evaluated in four different categories (cf. table I).

Assuming that the uncertainties be homogeneous for each category of test, we apply equations 5 and 7 and therefore need to estimate the values of σ_j . The estimates are given in table II (second column) and are based on the following considerations (items labelled according to the test category, cf. tables I and II – for a description of French grades look at table’s II caption):

p assuming $\sigma_p = 1$ amounts to saying that an error in grading by ± 3 units has a probability of occurrence $P < 0.003$ [9]. Translated in percentage $3\tilde{\sigma}_p = 0.15$, which is quite a large interval. Such a large error bar is introduced on the basis of the nature of the evaluation: different lab monitors give an estimate of the performance of each individual student – working in a small group (typically two or three) – by observing their work, discussing with the group and asking occasional questions. Each monitor is required to follow several groups (typically between four and six) and differences in evaluation among monitors, as well as fluctuations for a same monitor due to variable working conditions, are unavoidable.

l $\sigma_l = 0.5$ amounts to assigning ± 1.5 points to the uncertainty with probability $P > 0.997$ of the true grade falling within the interval. In percentage this amounts to $3\tilde{\sigma}_l = 0.075$. The variability in the evaluation is estimated to be lower than for p-tests due to the fact that labbooks, as a written document, can be more reliably evaluated. The estimated uncertainty could be smaller if all grading were done by one and the same person (not the case in this context). The size of σ_l is therefore chosen to reflect the added variability coming from an ensemble

TABLE II: The French University Grading System (FUGS) attributes the grades in $x/20$ (passing grade $x = 10$), where x represents the grade attributed to the test. The numerical estimates are given in absolute values, i.e., in FUGS units, but – in order to improve readability – are also repeated in percentage. The latter are identified by the corresponding quantities marked by a tilde \tilde{v} (v being the generic variable). The conversion gives rise to a result with an excess of digits for some grade categories (kept here to be consistent with the absolute estimates, used for the calculations).

Label	σ_j	σ_{G_j}	$\tilde{\sigma}_j$	$\tilde{\sigma}_{G_j}$
p	1	0.4	0.05	0.02
l	0.5	0.2	0.025	0.01
r	0.5	0.2	0.025	0.01
o	0.7	0.7	0.035	0.035

of graders.

r We assign the same error estimates to this test as those chosen for σ_l for the reasons exposed in the preceding point.

o This kind of test requires a closer look at its details. Being an oral examination – even though conducted by a panel of (at least) three examiners – it is somewhat more susceptible to fluctuations (in the questions, their evaluation, and in the student’s reactions). Therefore, we assign to it a value of $\sigma_o = 0.7$ which amounts to considering a full 95% confidence interval[10] ($\pm 2\sigma_o$) [1] to a spread of (approximately) three points. Translated into percentages, $\tilde{\sigma}_o = 0.035$ (i.e. $3\tilde{\sigma}_o = 0.105$). In other words, we expect the probability of a grade outside this interval to be below 5%.

The propagation of the *a priori* uncertainties on each individual grade for each kind of test follows equation 5 and produces the values of σ_{G_j} given in Table II. Notice that $M = 1$ for the oral test (label o), therefore no uncertainty improvement ensues for this grade.

Computing the propagation of the *a priori* uncertainty on the final grade, equation 7, we obtain

$$\sigma_G = 0.3 \text{ points}, \quad (8)$$

$$= 0.015 \text{ (in percent)} \quad (9)$$

which amounts to a $\Delta G \simeq 0.5$ points (or $\widetilde{\Delta G} \simeq 0.025$) with probability $P \simeq 0.9$ [11] of obtaining the actual grade within this interval.

We thus conclude that the *a priori* estimate for the grade each student receives in this course is ± 0.5 points (or 2.5%) with a confidence level of 90%.

V. DISCUSSION

Aside from the numerical result just obtained, equation 8, it is very instructive to look at the details of the

TABLE III: Individual contributions to the final grade uncertainty, estimated for each grade category according to equation 6 ($\frac{\partial G}{\partial G_j} = w_j$). We remark that only the last term in the table is significant.

	p	l	r	o
$w_j^2 \sigma_{G_j}^2$	0.0008	0.001	0.005	0.08

contributions which compose the final value σ_G . Table III provides the breakup of the various contributions, where we notice that the smallest one comes from the p component. We immediately recognize that the influence on the final uncertainty coming from the participation grade (p) is entirely negligible (by two orders of magnitude), in spite of its intrinsic variability and of the large *a priori* uncertainty we have consequently assigned to it. This results from the combined effect of the larger number of tests in this category ($M_p = 12$) and of the small weight assigned to this category ($w_p = 0.1$, cf. table I).

The relative contributions of the labbook (σ_l) and report (σ_r) uncertainties are different in table III due to their different number of samples ($M_l = 12$ and $M_r = 3$), which reduce by $M^{-\frac{1}{2}}$ their uncertainty (equation 5). Overall, however, even the weighted contribution coming from G_r is negligible – by more than one order of magnitude – when compared to that of the oral exam. We thus conclude that only the uncertainty on the latter matters in the determination of the uncertainty on the final grade, owing to the larger size of σ_o , the single event ($M_o = 1$), and especially its dominant weight ($w_o = 0.4$, table I).

One should not confuse the influence of each grade category on the final outcome with the dominance of the uncertainty of the oral test on the overall uncertainty. Each grade category contributes, proportionally to its weight, to the course grade, but the confidence interval is determined exclusively by the oral test, all other forms of grading providing a much more “accurate” evaluation.

This result has the following implications:

- given the very sizeable difference in error contribution (table III) modifying the estimates of the *a priori* uncertainties which we have assigned to the various kind of tests (except for the oral test) will not influence the size of the uncertainty. Thus, except for the oral test, we realize *a posteriori* that our careful estimates in the preceding section do not hold any relevance;
- given that the only dependence of the estimated *a priori* final uncertainty has a linear dependence onto the estimated error assigned to the oral test, we know that modifying the latter linearly translates onto the reliability of the global grade (weighted by w_o);
- there is no need to worry about the reliability of the grades for the first three kinds of tests, i.e. about

the variability originating from the involvement of several monitors in the various grading steps.

The last point is important for instructors who may worry that, in particular for the participation grade, the intrinsic variability due to multiple evaluators, and the ensuing point spread, may distort the reliability of the course’s global grade. This also means that one can confidently introduce different measures of evaluation – in particular some which give the benefit of an immediate return to the students, such as the participation grade – without risking a substantial impact onto the final grade.

One final point: the *a priori* estimate of the uncertainty on the final grade also gives a measure of the precision with which the latter can be given. In the specific case of the example used, a good discretization is ΔG , i.e., using a scale in points with integer and half-integer values (or 2.5% in relative precision). This can be used to explain to students what is a reasonable scale in grade spacing. Of course, the actual value will depend on the structure of the course and on the number of test categories (and of test number in each category).

VI. CONCLUSIONS

The simple, but striking, application of the *APA* to course grades illustrates quite effectively the intrinsic nature of this kind of analysis and its main features. It clearly shows the technique’s importance in all those situations where measurements cannot be repeated to obtain *a posteriori* error estimates, and the power of its predictions. At the same time, the analysis has shown the need for a careful assessment of the individual error sources to be assigned to the primary, *measured* quantities, and the futility of part of the work, rendered irrelevant by the intrinsic structure of the analyzed quantities (composition of the grade and of its uncertainty). Students are often taken aback both by the fact that an aspect of their curriculum can be analyzed in detail with techniques *seemingly exclusively devised* for lab experiment analysis, and by the information which can be gathered by this analysis. This example should also serve as an encouragement for testing the application of data treatment techniques to real-life everyday’s problems.

As a bonus, we have shown that instructors may gather precious information on the uncertainty which affects their grades and on the confidence level of each grade component, while gaining some freedom in experimenting with innovative ways of introducing partial grades. The latter can be beneficial to giving students early and welcome feedback, while ensuring that the reliability of the course grade is not affected by evaluation components which are more prone to larger fluctuations. This quantitative analysis, though partly subjective (in the assignment of the *a priori* error components) may also be helpful in arguing in favour of the introduction of complementary grade parts in the discussions with skeptical colleagues or Department Directors.

I am grateful to all the students (in excess of five hundred) who, having taken this course over the last decade,

have stimulated the development of new ideas and examples.

-
- [1] Taylor J R 1982 *An introduction to uncertainty analysis: the study of uncertainties in physical measurements* (University Science Books, Mill Valley, CA)
 - [2] Dwight H B 1961 *Tables of Integrals and Other Mathematical Data*, fourth edition (Macmillan, New York)
 - [3] Not too many students in a class will be confronted with the needs of a true *APA* in the context of their future careers.
 - [4] It is also useful to check whether the errors obtained are consistent with those expected, this enabling an independent test of the error size and, possibly, the detection of mistakes (or systematic errors). This is often beyond the reach of a student lab, though.
 - [5] A good example within the physical sciences can be given when considering working on existing datasets – e.g., old climatologic or meteorological data, where only measurements without uncertainties may have been registered.
 - [6] This is the structure of grades of the laboratory course which I use as an example for the students.
 - [7] This can be useful when one test from an ensemble proves more difficult to grade, or exhibits larger fluctuations in students' performance (e.g., due to a harder problem set).
 - [8] The nature of the problem is the same if different uncertainties are attributed to the individual test, as previously mentioned.
 - [9] Gaussian-distributed errors are assumed here [1].
 - [10] It is interesting here to use a different way of evaluating the confidence interval. Indeed, contrary to common practice in Physics, in most applied sciences – e.g. Risk Assessment, Geosciences, etc. – the standard error bar associated with any given quantitative estimate is $\pm 2\sigma$, rather than $\pm\sigma$ by virtue of its larger confidence level (95%), more useful for practical purposes. It is therefore instructive to use this example to present this aspect of uncertainty estimates to physics students.
 - [11] From [2] we obtain that $P(x) = 0.9$ when $x \simeq 1.645$. Thus, computing $\Delta G = 1.645 \times \sigma_G$ (equation 8) we obtain $\Delta G \simeq 0.4935 \approx 0.5$, as in the text.